# A Meta-Analysis Approach for Feature Selection in Network Traffic Research

Daniel C. Ferreira, Félix Iglesias Vázquez, Gernot Vormayr,
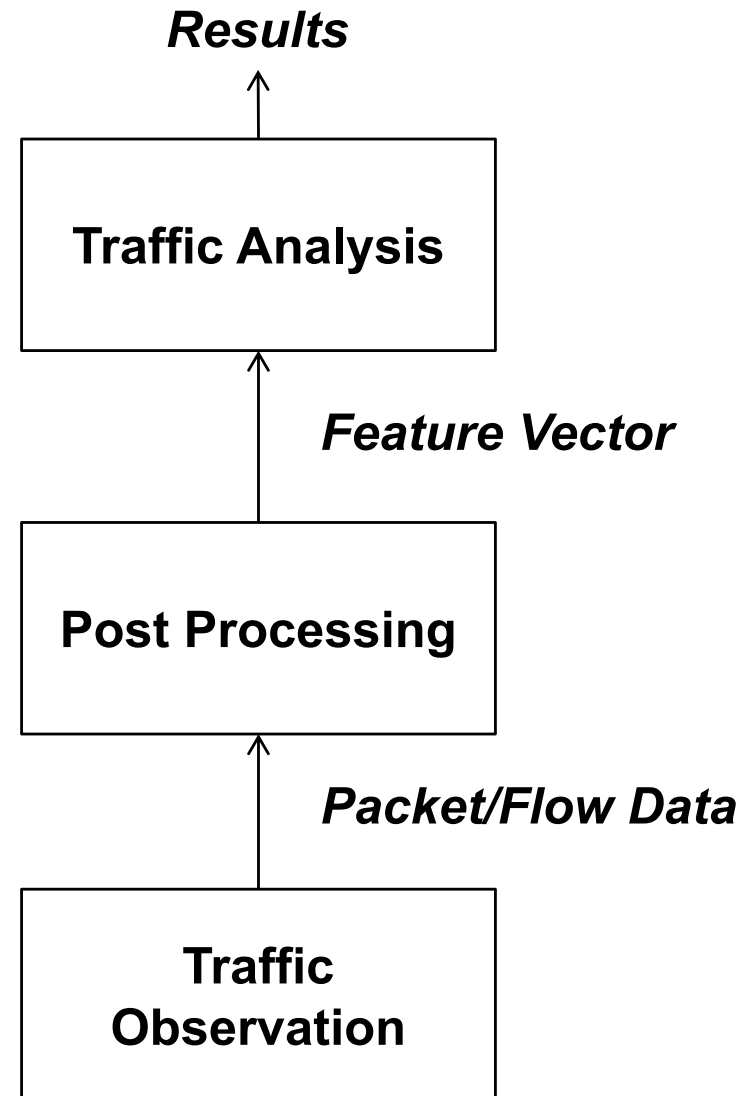Maximilian Bachl, Tanja Zseby
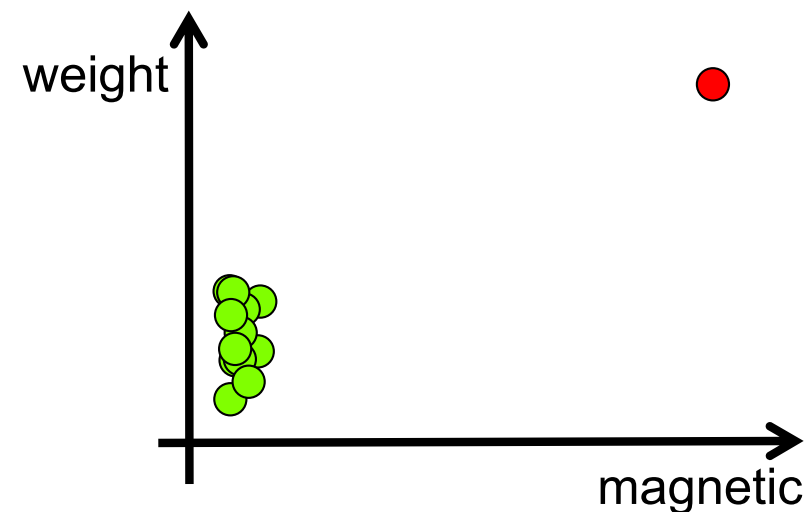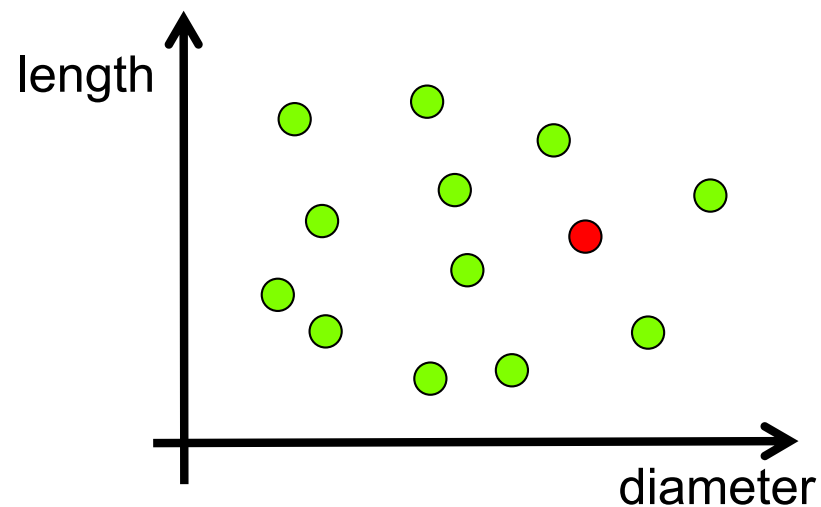
Institute of Telecommunications
TU Wien

institute of
telecommunications

TU WIEN
TECHNISCHE
UNIVERSITÄT
WIEN
Vienna | Austria

# Network Traffic Analysis

**Feature Selection:**

Select most suitable features

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ ... \end{bmatrix}$$



*Results*

**Traffic Analysis**

*Feature Vector*

**Post Processing**

*Packet/Flow Data*

**Traffic Observation**

# Well-chosen Features ➔ Simplified Analysis



length

diameter

weight

magnetic

institute of
telecommunications

# Agree to Disagree



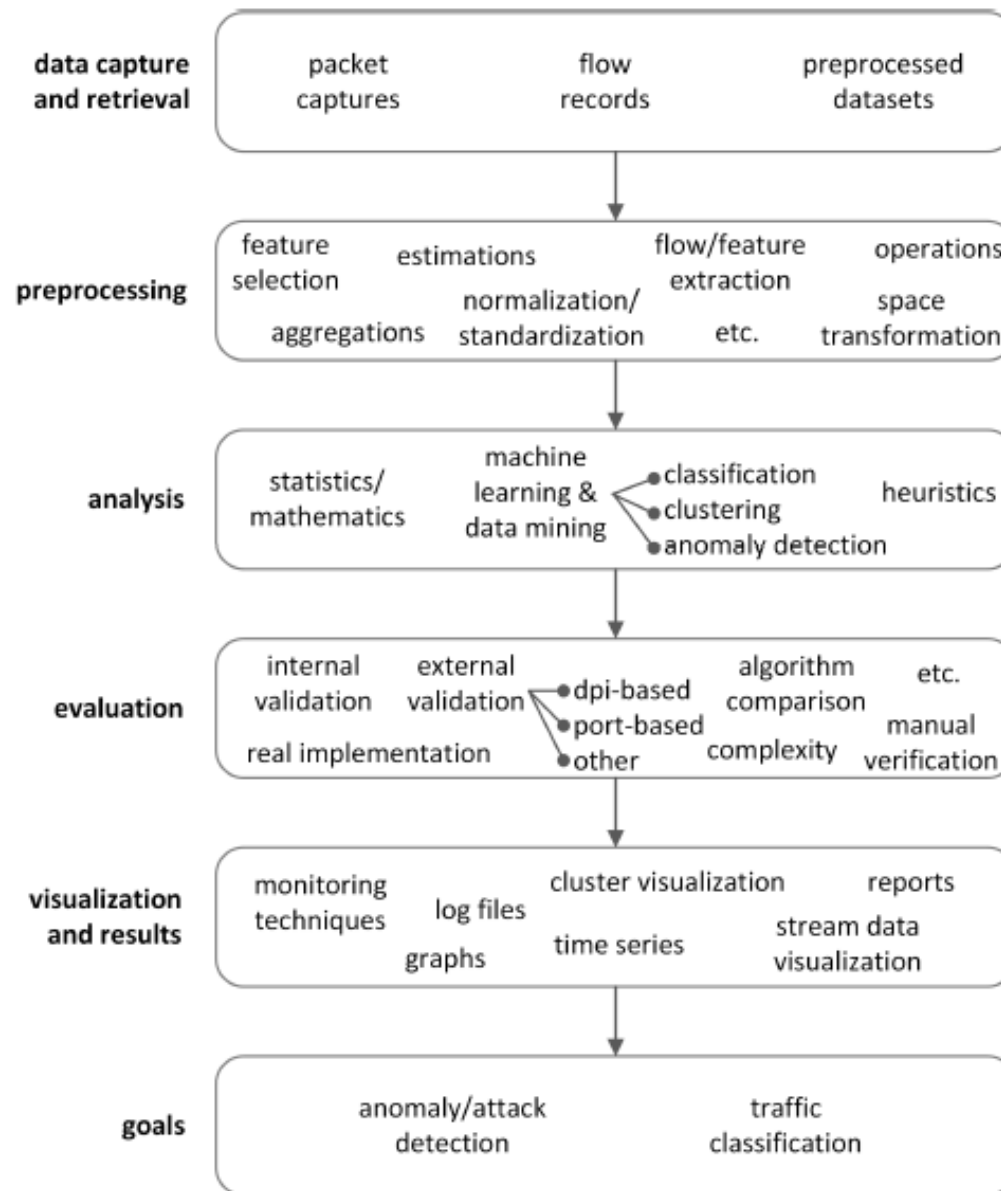Source: Iglesias, Zseby: "*Analysis of network traffic features for anomaly detection*"; Machine Learning, **101** (2015), 1; 59 - 84.

# Why a Meta Analysis?

- Meta-Analysis common in other disciplines
  - Structures the state of art
  - Combines existing results
  - Identifies agreements/disagreements in the community
  - Provides basis for gap analysis
- Provides information about
  - Availability of data and tools
  - Parameter settings
  - Validation Methods
  - Terminology and notation

➔ Supports reproducibility and comparability

# Data Structure

# Example: Features

- Base features
- Operations on base features
- Flow keys

*Standard IPFIX Information Element*

```
"features": [
    {"log": ["octetTotalCount"]},
    {"log": [{"divide": ["octetTotalCount", "_activeForSeconds"]}]},
    {"maximum": ["_interPacketTimeMicroseconds"]},
    {"minimum": ["_interPacketTimeMicroseconds"]},
],
"key_features": [
    "sourceIPv4Address",
    "destinationIPv4Address",
    "protocolIdentifier"
]
```

*Non-IPFIX feature*

# Example: Data Set

```
"data": {
  "datasets": [
    {
      "dataset_name": "mawi-2015",
      "availability": "public",          ←——————— Dataset available
      "format": "packet",
      "types": "ip",
      "generation": "captured",
      "generation_year": 2015,
      "covered_period": "minutes",
      "details": ["raw","no_payload"],
      "subsets": ["01-01-2015","15-04-2015","31-07-2015"]
    },
```

institute of
telecommunications
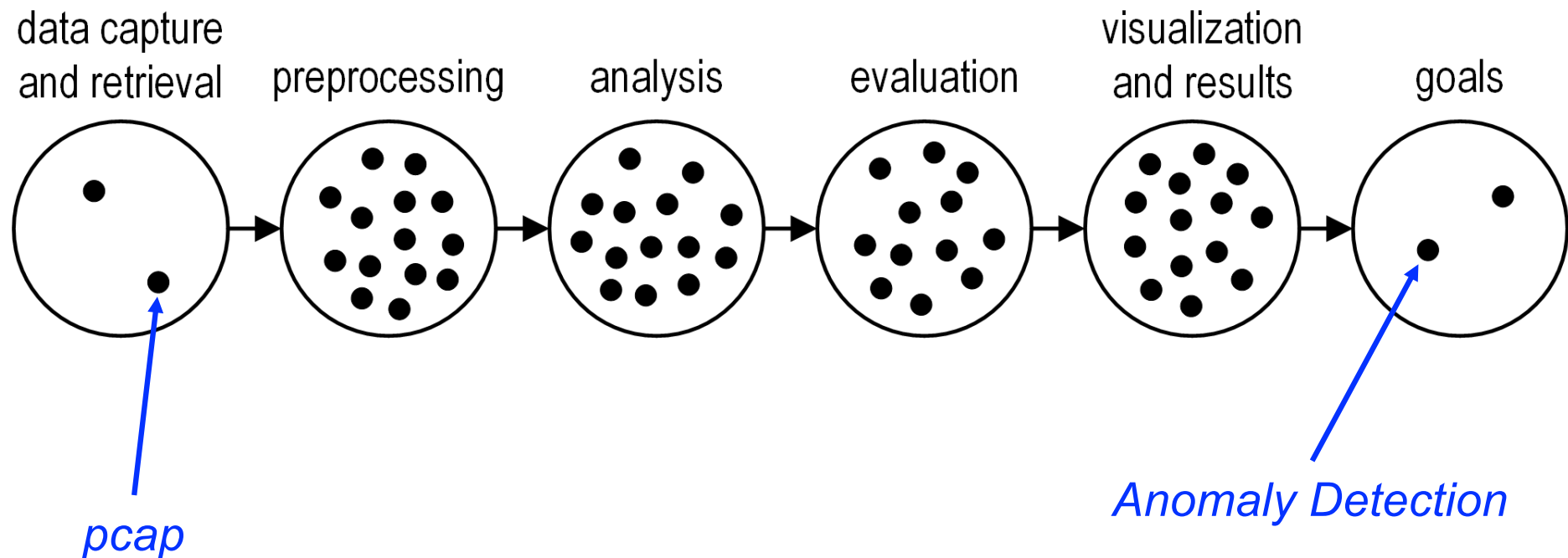
# Example: Algorithms

```
"algorithms": [
    {
        "name": "fuzzy clustering",
        "subname": "gustafson-kessel",
        "learning": "unsupervised",
        "role": "main",
        "type": "clustering",
        "metric/decision_criteria": "mahalanobis",
        "tools": [
            {
                "tool": "matlab_fuzzyclusteringtoolbox",
                "detail": "none",
                "availability": "public"          ← Tool available
            }
        ],
        "source": "referenced",                   ← Link to tools provided
        "parameters_provided": false              ← Parameters not
    },                                               provided
```

# Initial Results

- 71 Papers from years 2005 to 2017

## Analysis Chain

# Initial Results

- Flow Definitions
  - 64.6% of papers that define a flow-key use classical 5-tuple {sIP, dIP, sPort, dPort, Protocol}
  - 70.8% use bi-directional flows
  - 83.1% use flow-based features
- Data Sets
  - 46.5% use at least one public data set
- Most Common Features
  - Number of papers that use a specific base feature
  - Number of papers weighted with their citations $\log_{10}(\text{citations})$

institute of
telecommunications

# Most Recurrent Base Features

| Features (recurrences) | score[1] | Features (citations) | score[2] |
|---|---|---|---|
| octetTotalCount | 5.8 | octetTotalCount | 4.6 |
| packetTotalCount | 3.9 | ipTotalLength | 3.9 |
| flowDurationMilliseconds | 3.1 | destinationTransportPort | 3.5 |
| ipTotalLength | 2.7 | sourceTransportPort | 3.0 |
| destinationTransportPort | 2.5 | flowDurationMilliseconds | 2.6 |
| destinationIPv4Address | 2.4 | packetTotalCount | 2.3 |
| sourceIPv4Address | 2.3 | destinationIPv4Address | 2.3 |
| sourceTransportPort | 2.0 | sourceIPv4Address | 2.3 |
| protocolIdentifier | 2.0 | protocolIdentifier | 2.2 |
| _interPacketTime$\mu s$ | 2.0 | _server_to_client | 2.2 |
| _server_to_client | 1.5 | _client_to_server | 2.2 |
| _client_to_server | 1.5 | _interPacketTime$\mu s$ | 1.8 |

# Summary

- Meta Analysis for Network Traffic Analysis
  - Supports comparability and reproducibility
  - Focus on feature selection, but much more information collected
- JSON files
  - Structured, searchable state of art
  - Fast extraction of relevant information from papers
- Initial results
  - Most common features
  - Flow definitions
  - Usage of public data sets
- → Data allows for many further analysis opportunities

# Discussion

- Manual data curation ➨ Errors
  - Involve authors (check and correct)
- Analysis just shows "preferred" features, methods
  - ➨ not necessarily the best!
- Incentives to fill data base

  - Conferences can require to add accepted papers
  - Students can add data when exploring state of art
  - Searchable data base may increase citations for papers included

- All data, documentation, paper data base available at: *www.cn.tuwien.ac.at/meta*

# Thank you!

Contact: tanja.zseby@tuwien.ac.at